



# GDS : Grid Data Service

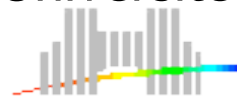
## Un service de gestion de données pour les grilles de calcul

Gabriel Antoniu<sup>1</sup>, Eddy Caron<sup>2</sup>,  
Frédéric Desprez<sup>2</sup>, Pierre Sens<sup>3</sup>

<sup>1</sup>IRISA, Rennes

<sup>2</sup>LIP, ENS Lyon

<sup>3</sup>LIP6, Université Paris 6



*Journées PaRISTIC, Bordeaux, 21-23 novembre 2005*



# Le projet GDS

---

- Projet de l'ACI Masses de Données (2003)
- Objectif :

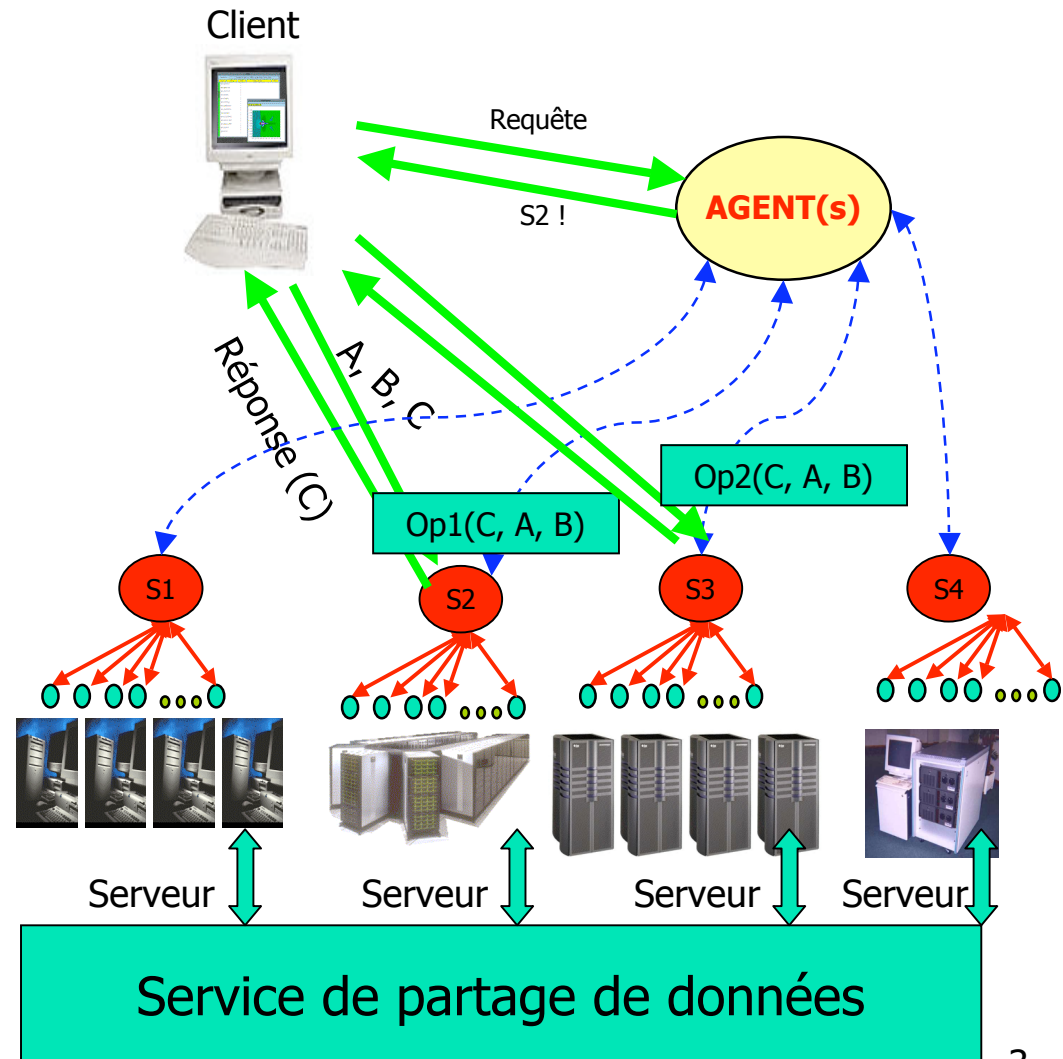
*Spécifier, réaliser et valider expérimentalement un service de partage de données pour la grille, adapté au calcul scientifique*

- Partenaires
  - Projet GRAAL (LIP, Lyon)
  - Projet PARIS (IRISA, Rennes)
  - Projet REGAL (LIP6, Paris)

# GDS : objectifs

- Partager des données entre plusieurs composants applicatifs distribués sur une grille

- Transparence de la localisation des données
- Stockage persistant
- Tolérance aux fautes
- Gestion de la cohérence des copies en environnement dynamique
- Extensibilité à l'échelle d'une grille : des milliers de nœuds

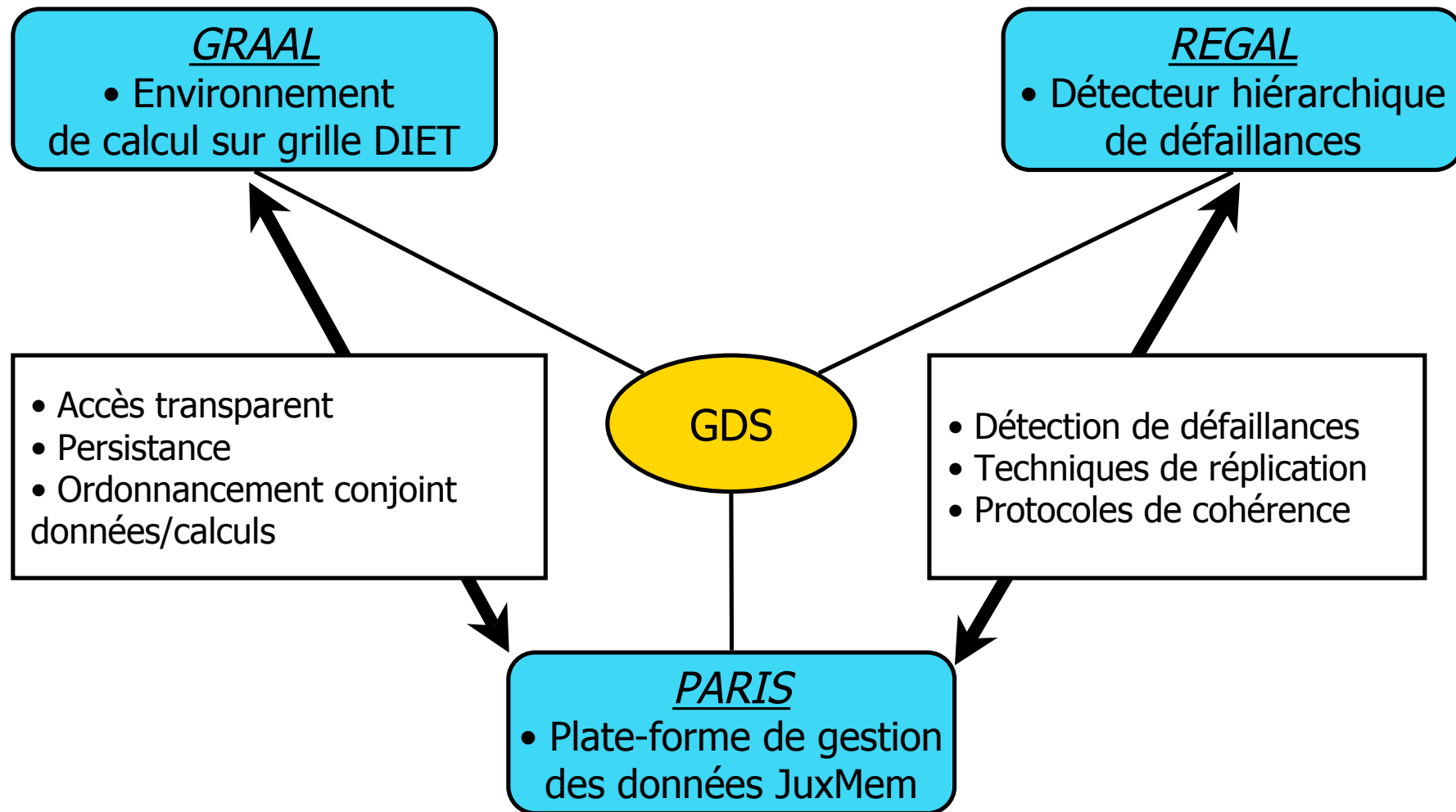


# Idée : une approche hybride

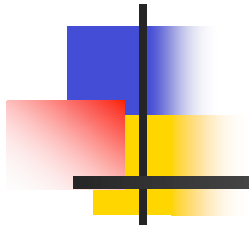
- **Systemes MVP (DSM)** : cohérence et accès transparent
- **Systemes P2P** : extensibilité et volatilité

	<b>MVP</b>	<b>Service de partage pour la grille</b>	<b>P2P</b>
<b>Plate-forme physique</b>	Grappe	Fédération de grappes	Internet
<b>Echelle</b>	$10^1$ - $10^2$	$10^3$ - $10^4$	$10^5$ - $10^6$
<b>Volatilité</b>	Nulle	Moyenne	Haute
<b>Contrôle</b>	Fort	Moyen	Faible
<b>Applications</b>	Calcul scientifique	Calcul scientifique et stockage	Partage et stockage de fichiers
<b>Type de données</b>	Modifiables	Modifiables	Non modifiables
<b>Tolérance aux fautes</b>	Rare	?	Critique
<b>Cohérence</b>	Critique	?	Rare

# GDS : interaction des partenaires

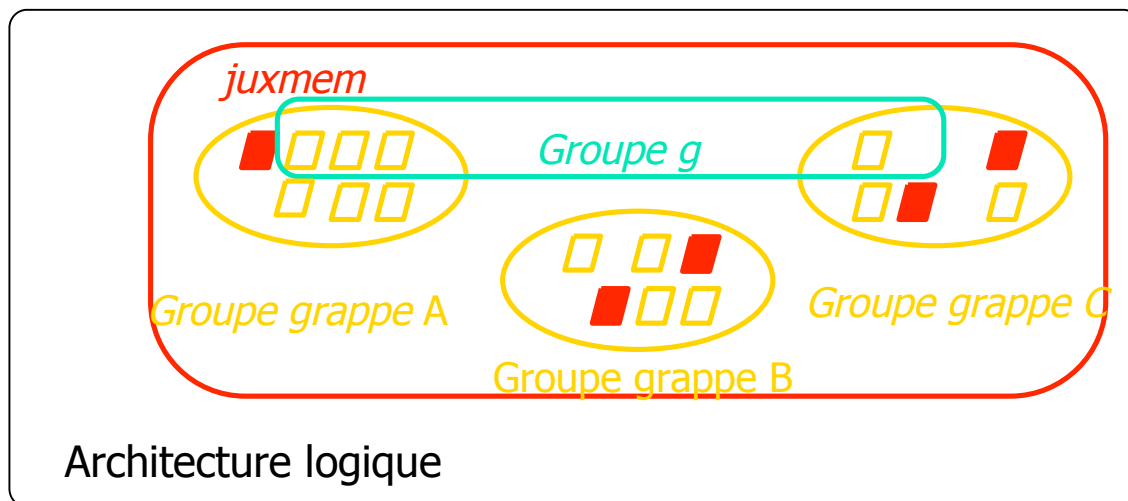


# JuxMem



# JuxMem

- Une architecture pair-à-pair de service de partage de données en mémoire
- Mécanisme de persistance
- Localisation transparente des données
- Mécanisme de réplication
- Mécanismes de cohérence



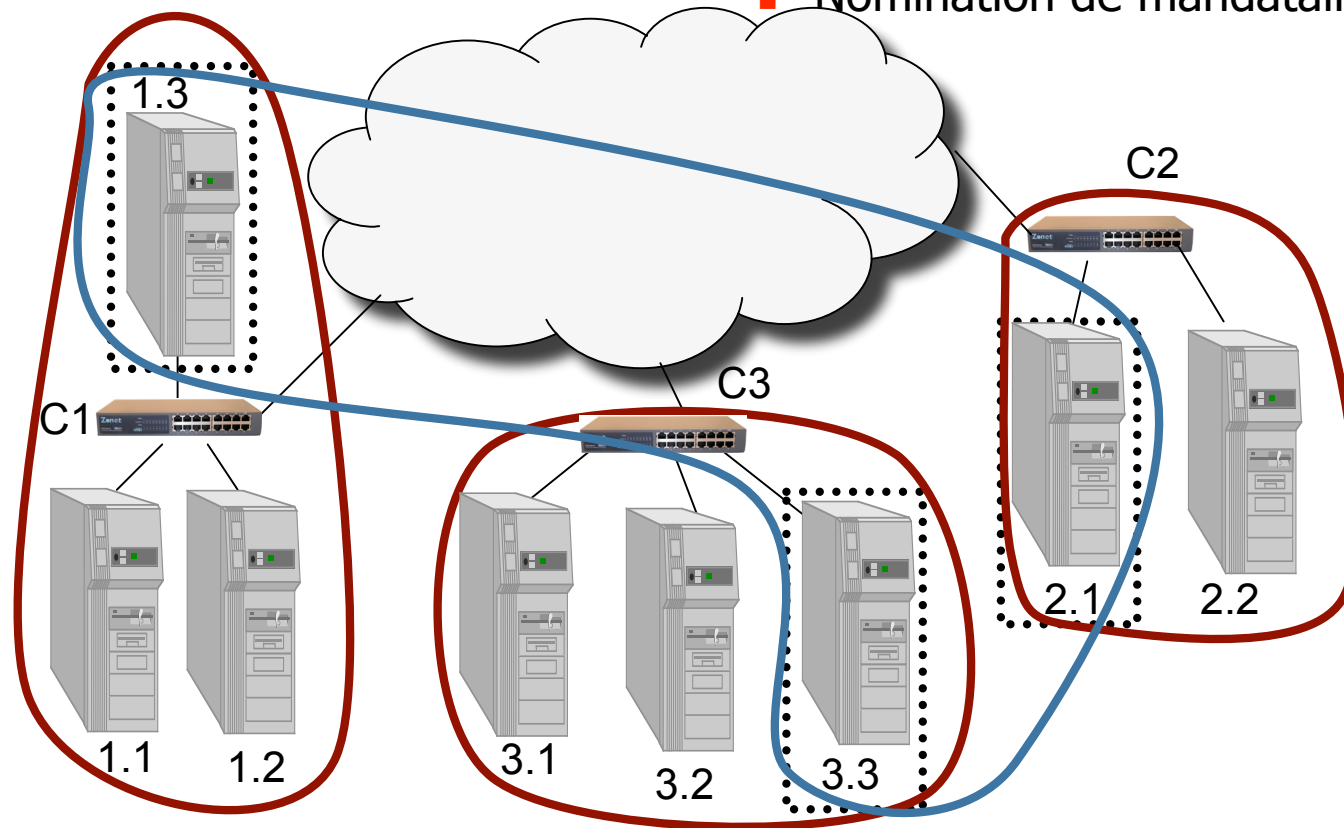
- Plateforme de programmation d'application P2P
  - Ensemble de protocoles
- Un pair
  - Identifiant unique (ID)
  - Adresse indépendante de la localisation physique
  - Plusieurs protocoles d'accès réseau (TCP, HTTP, etc)

Project  
**JXTA**

<http://www.irisa.fr/paris/Juxmem>

# Tolérance aux fautes : organisation hiérarchique

- 1 groupe local / cluster
- 1 groupe global
- Composition du groupe global
  - Représentant dans chaque groupe local = mandataire
  - Nomination de mandataire

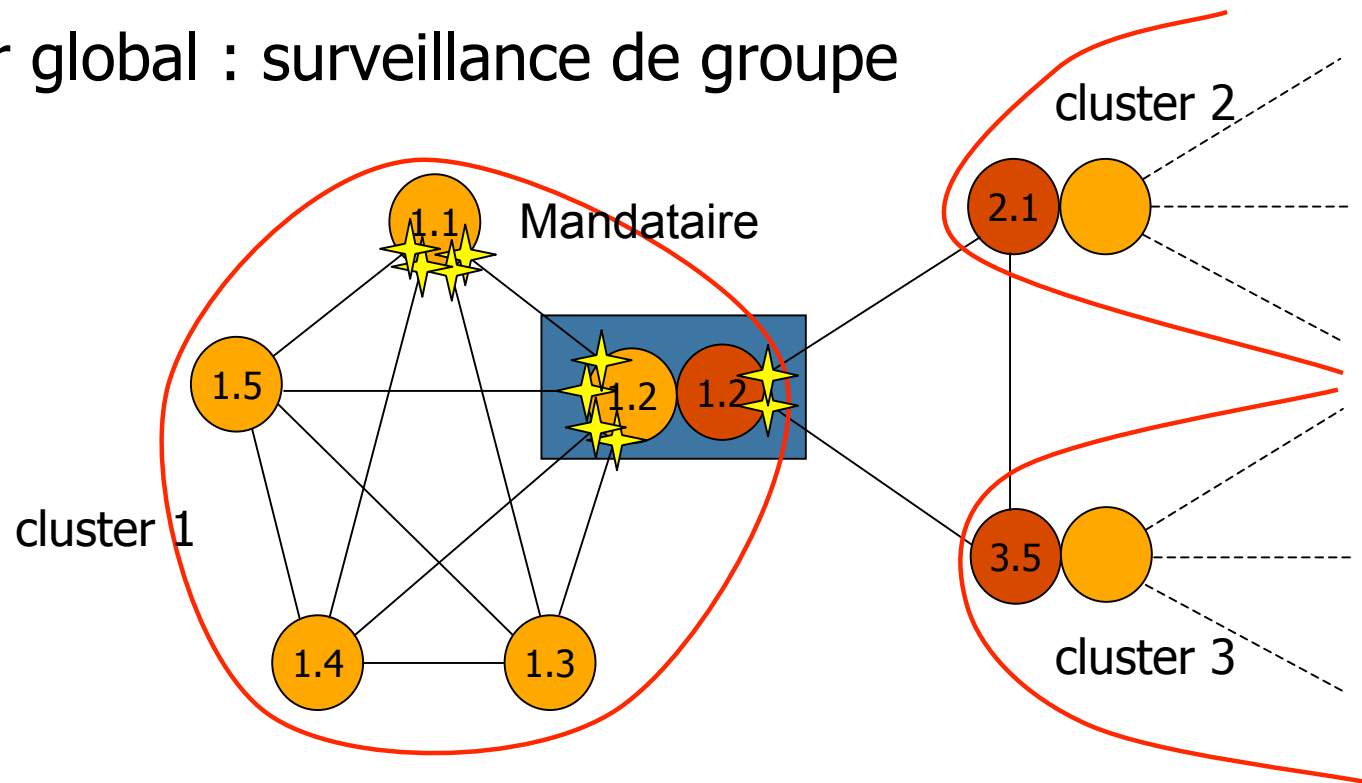


Organisation  
hiérarchique



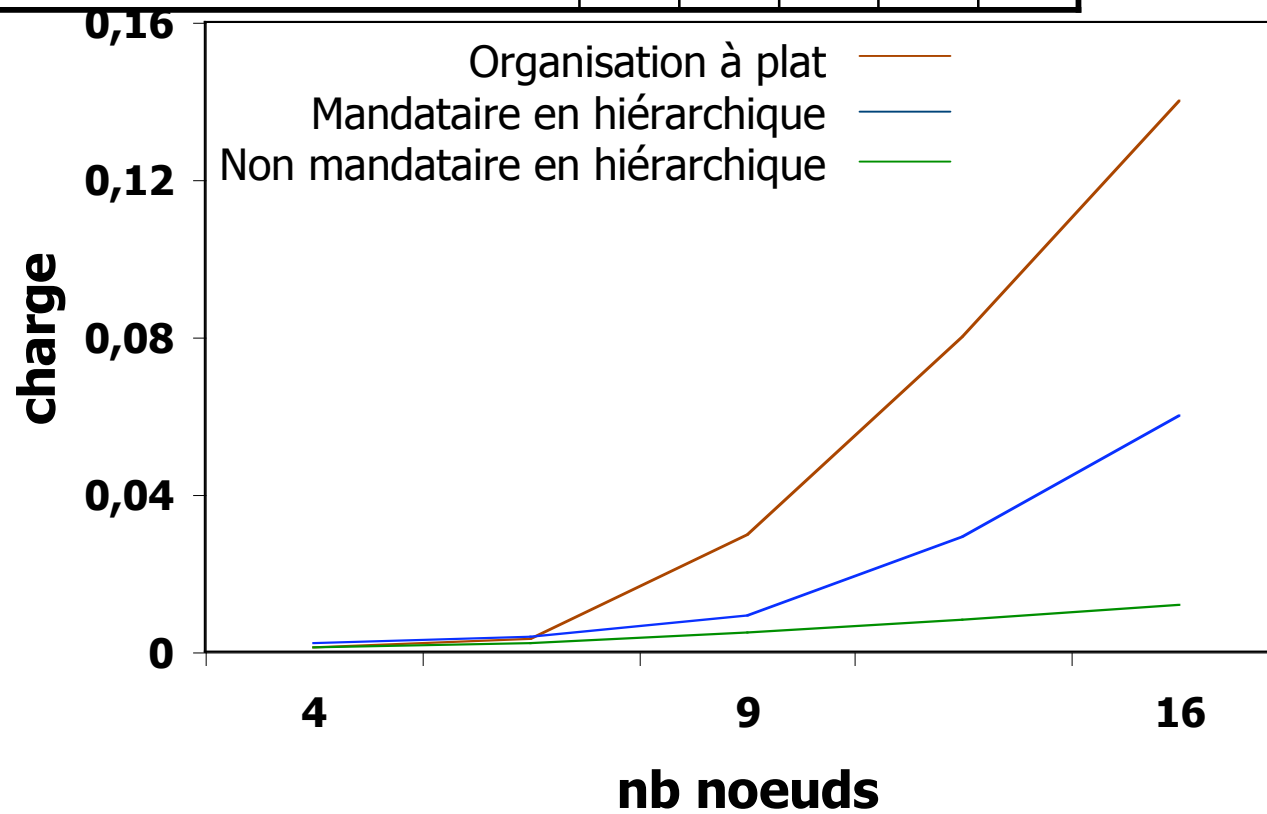
# Rôle du mandataire

- Permet la connexion du groupe local avec le reste du système
- Détecteur local : surveillance de nœud
- Détecteur global : surveillance de groupe

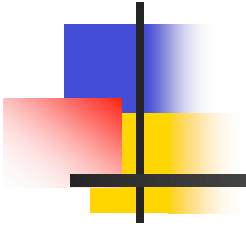


# Performances

Nb de nœuds	4	6	9	12	16
Nb de nœuds par groupe local	2	3	3	4	4
Nb de groupes locaux	2	2	3	3	4



Organisation  
hiérarchique



DIET

---



# DIET

---

- Our goals

- To develop a toolbox for the deployment of environments using the Application Service Provider (ASP) paradigm with different applications
- Use as much as possible public domain and standard software
- To obtain a high performance and scalable environment
- Implement and validate our more theoretical results
  - Scheduling for heterogeneous platforms, data (re)distribution and replication, performance evaluation, algorithmic for heterogeneous and distributed platforms, ...

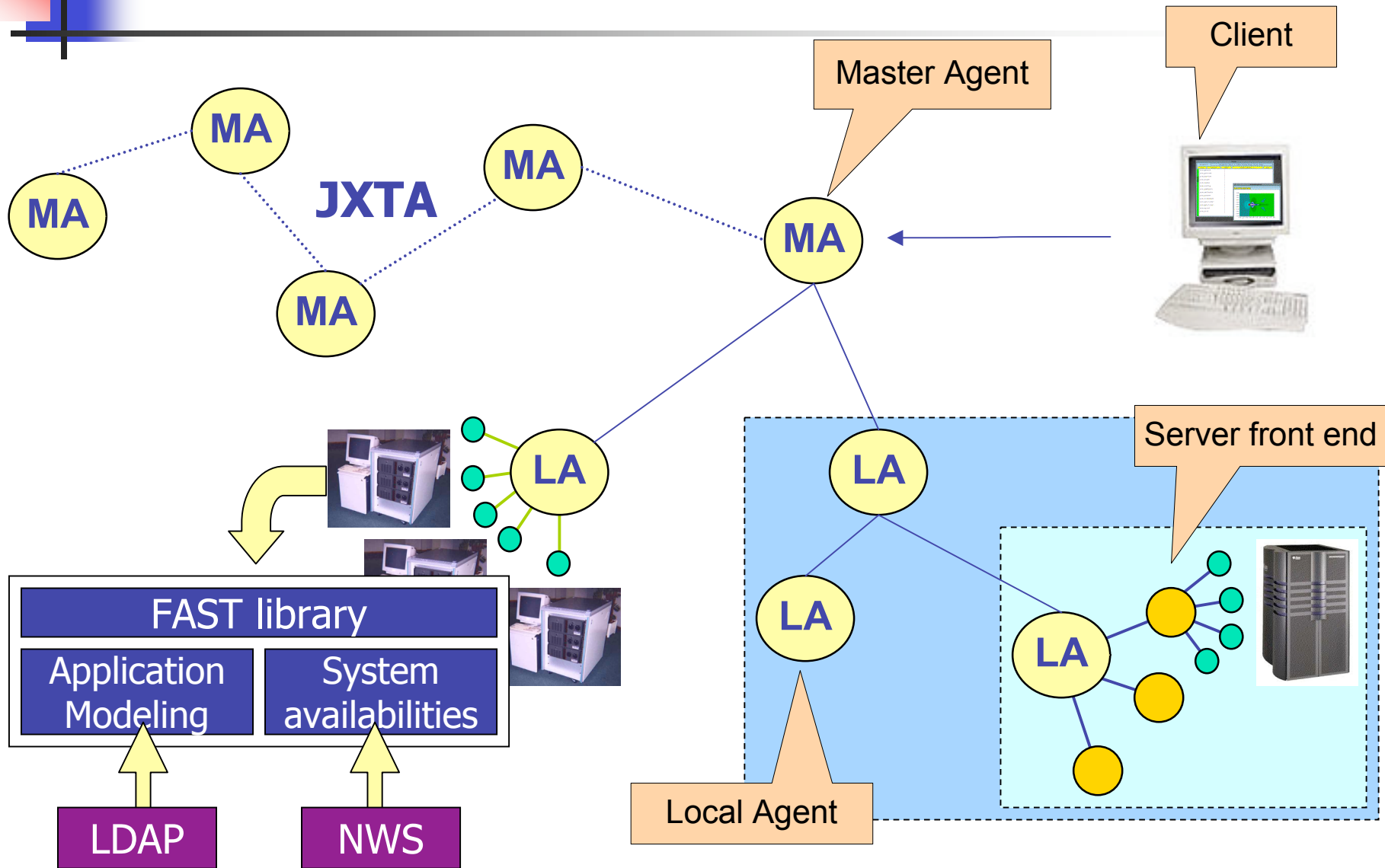
- Based on CORBA, NWS, LDAP, and our own software developments

- FAST for performance evaluation,
- LogMgr for monitoring,
- VizDIET for the visualization,
- GoDIET for the deployment

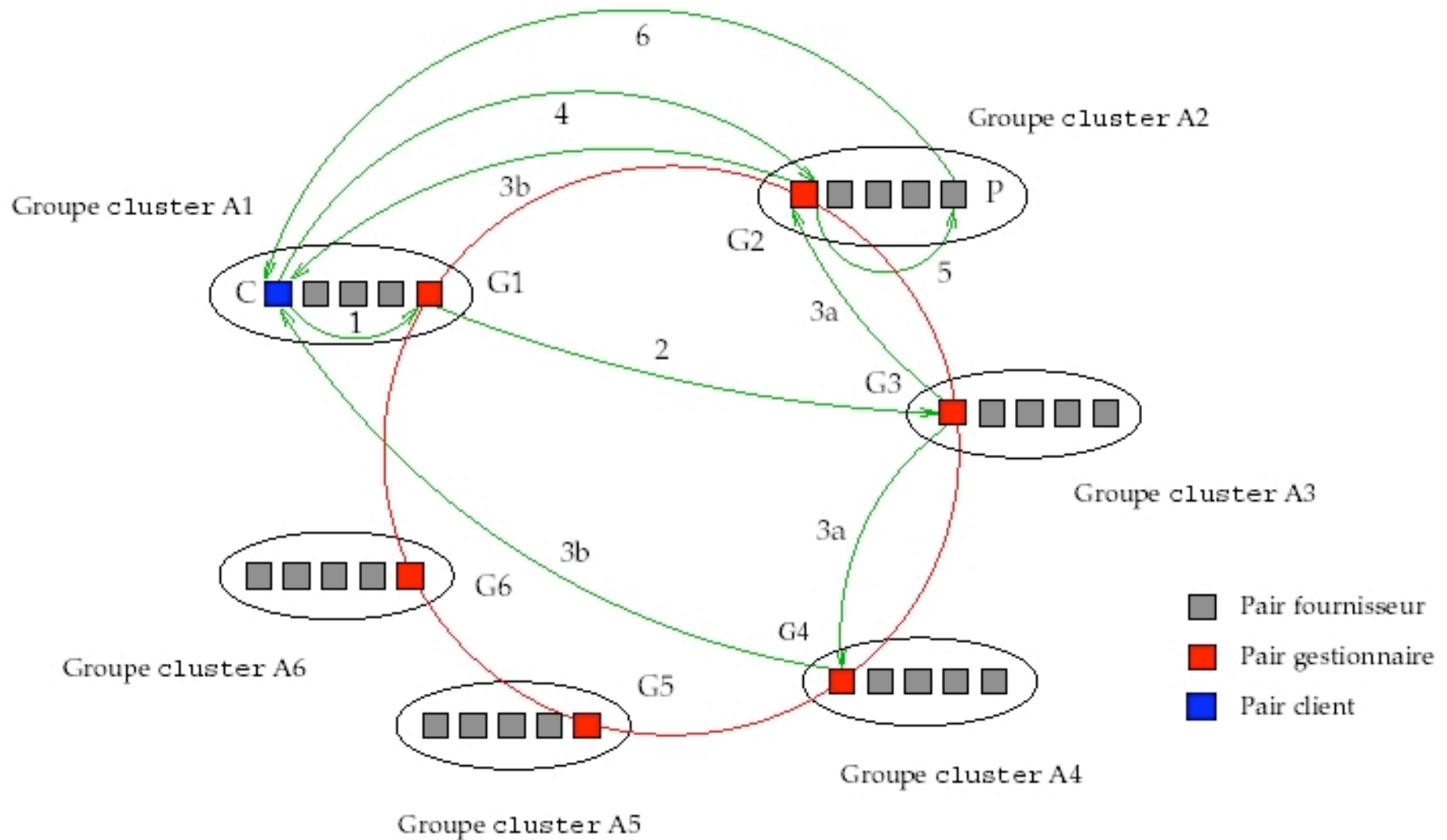
- Several applications in different fields (simulation, bioinformatic, ...)

- Release 2.0 available on the web
- ACI Grid ASP, RNTL GASP

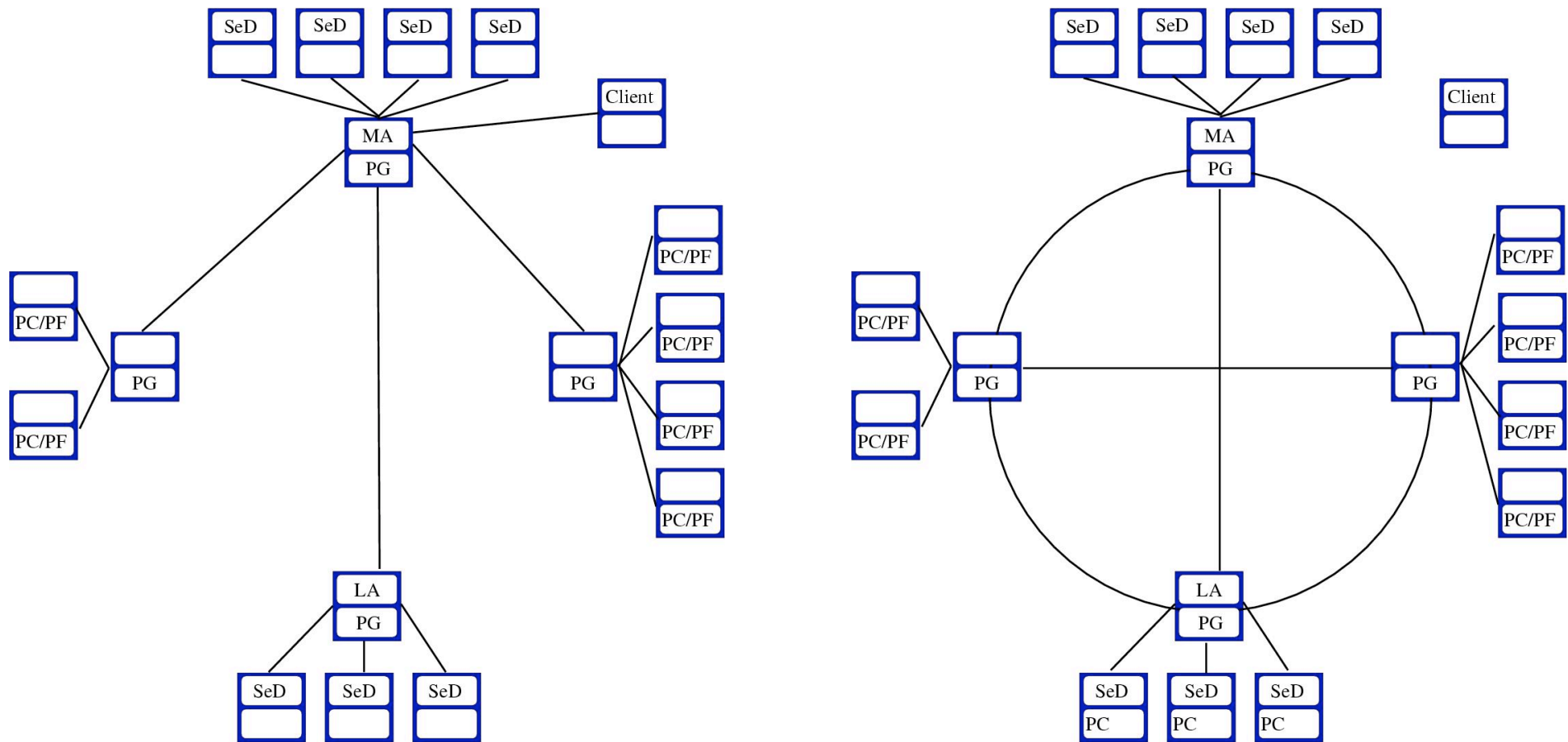
# Architecture DIET



# Architecture JuxMem

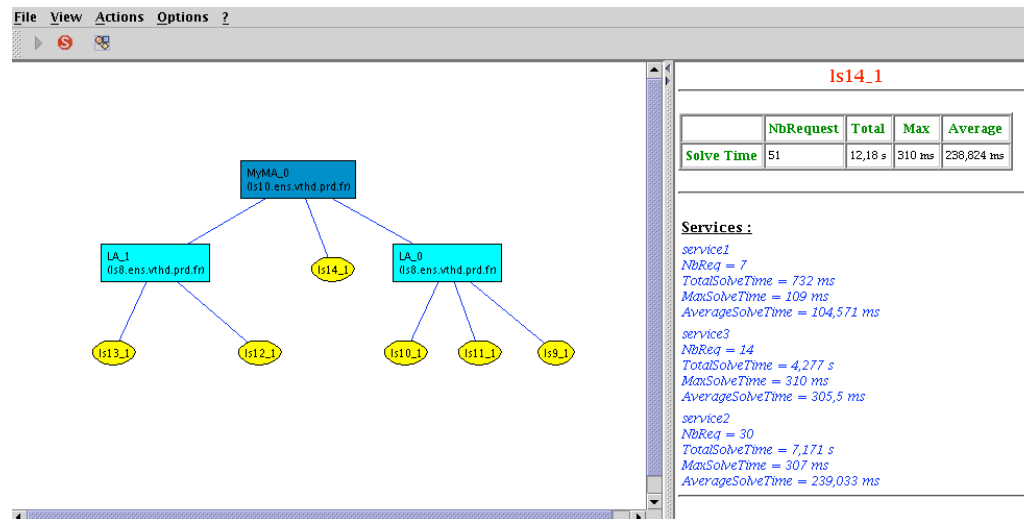
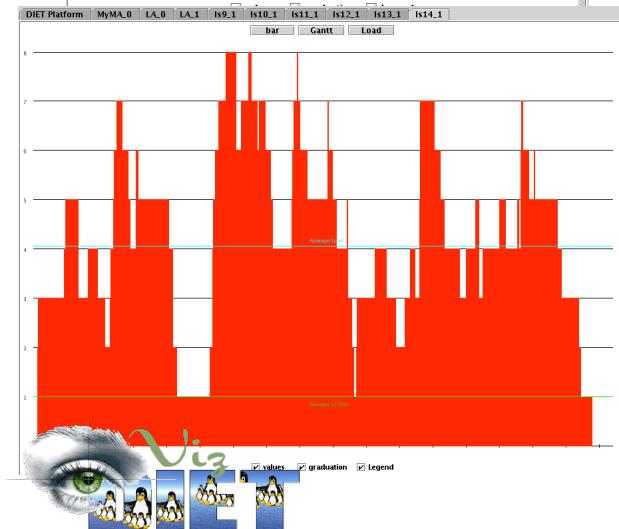
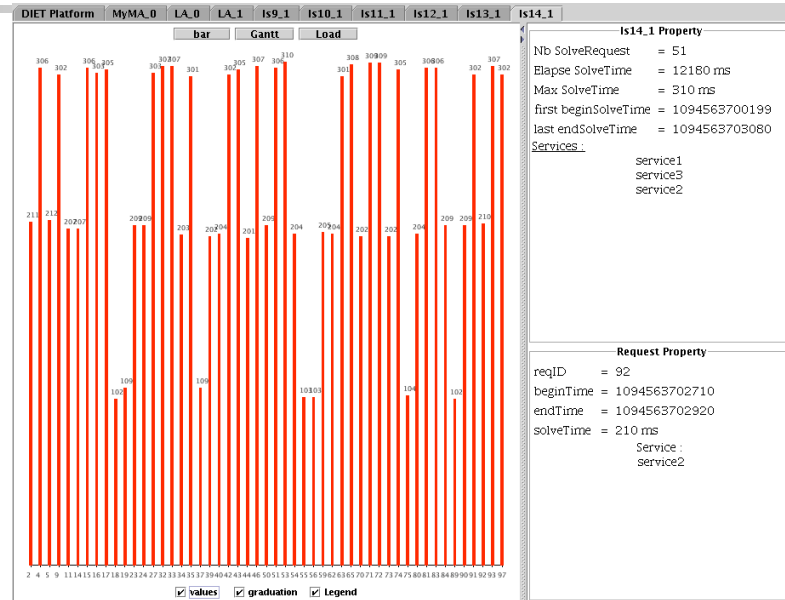
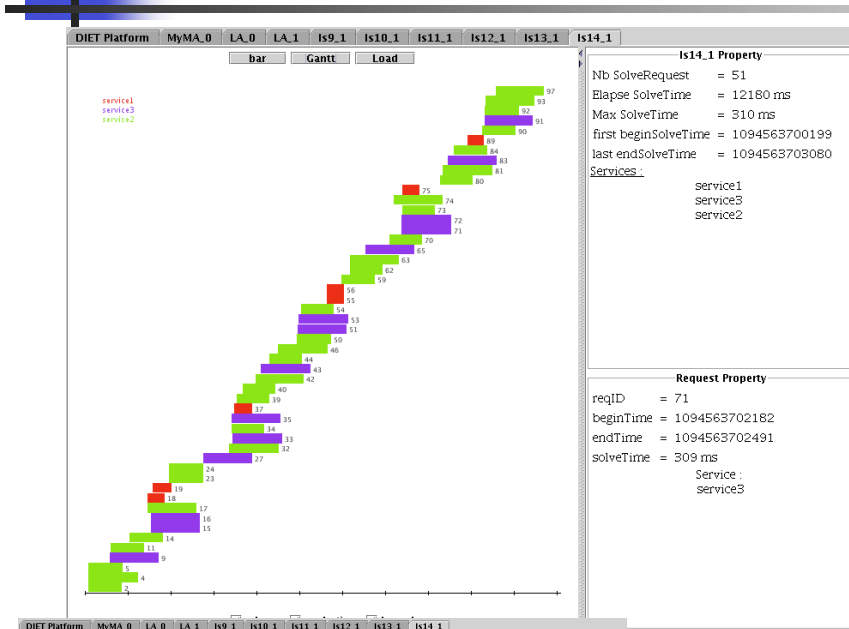


# Cohabitation DIET / JuxMem



- Données répliquées sur les grappes de la grille
- Implémentation multi-protocoles (réplication, cohérence)

# Screenshot





# DIET Dashboard

- Current works : We designed a new tool to merge GoDIET and VizDIET
- Resource visualization before GoDIET deployment
- Interactive GoDIET
- Only one tool to manage all DIET platform.
- DIET dashboard is designed from BuildDIET project



File Log Construct Utils View Help

Resources

- Compute resources
  - IsCluster
    - Is1
    - Is2
    - Is3
    - Is4
    - Is5
    - Is6
    - Is7
    - Is8
    - Is9
- cristalCluster
  - clus-103
  - clus-104
  - clus-105
  - clus-106
  - clus-107
  - clus-108
- Storage resources
  - IsDisk
  - cristalDisk

MA\_0\_0  
MA\_0\_0-clus-108LA\_0\_0-LA\_2\_0  
clus-105\_1  
LA\_2\_0  
LA\_2\_0-SetD\_0-SetD\_1-SetD\_2-SetD\_6\_0  
SeD\_7\_1  
SeD\_8\_2  
SeD\_6\_0

MyMA\_0  
(ls10.ens.vthd.prd.fr)

LA\_0  
(paraci60)

LA\_1  
(ls12.ens.vthd.prd.fr)

paraci12\_1  
paraci13\_1  
paraci14\_1  
paraci15\_1  
paraci16\_1  
paraci17\_1  
paraci18\_1  
paraci20\_1

paraci11\_1

	NbRequest	Total	Max	Average
Solve Time	8	2,552 s	341 ms	319 ms

Services :  
service3

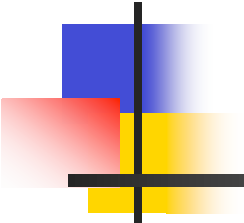
TotalSolveTime = 2,552 s  
MaxSolveTime = 341 ms  
AverageSolveTime = 319 ms

Persistent Data :

Local scratch directory ready.  
/tmp/run\_05mai22\_1923

\*\* Launching element OmniNames on Is2  
Writing config file omniORB4.cfg  
Staging file omniORB4.cfg to IsDisk  
Executing element OmniNames on resource Is2

Log Deploy External program



# Les équipes – avancement

---



# Projet GRAAL

INRIA Rhône-Alpes - CNRS – ENS Lyon

- Domaine de recherche :
  - Algorithmique et ordonnancement pour plates-formes hétérogènes distribuées
  
- Compétences de GRAAL (relatives à GDS) :
  - Environnements de calcul sur grille de type Network Enabled Servers (plate-forme DIET)
  - Ordonnancement des calculs
  - Algorithmes et outils de déploiement hiérarchique
  
- Personnels impliqués directement dans GDS :

<b>Permanents</b>	<b>Doctorants</b>	<b>Ingénieurs</b>
<ul style="list-style-type: none"><li>■ Eddy Caron (MDC ENS Lyon)</li><li>■ Frédéric Desprez (DR INRIA)</li></ul>	<ul style="list-style-type: none"><li>■ Pushpinder-Kaur Chouhan (MESR)</li><li>■ Cédric Tedeschi</li></ul>	<ul style="list-style-type: none"><li>■ Raphael Bolze (INRIA)</li><li>■ Holly Dail (INRIA)</li></ul>



# Projet PARIS

INRIA Rennes - CNRS - Université de Rennes 1

INSA de Rennes - ENS Cachan / Bretagne

- **Domaine de recherche :**
  - Modèles de programmation et supports exécutifs pour les systèmes parallèles et distribués (grappes, grilles, P2P)
  - Applications : simulations numériques distribuées
  - Gestion de données pour grilles de calcul et systèmes P2P
  
- **Compétences de PARIS (relatives à GDS) :**
  - Partage de données à grande échelle (plate-forme JuxMem, basée sur la plate-forme JXTA de Sun Microsystems)
  - Protocoles de cohérence des données
  
- **Personnels impliqués directement dans GDS :**

<b>Permanents</b>	<b>Doctorants</b>	<b>Stagiaires</b>
<ul style="list-style-type: none"><li>■ Gabriel Antoniu (CR INRIA)</li><li>■ Luc Bougé (Prof. ENS Cachan, Antenne de Bretagne)</li></ul>	<ul style="list-style-type: none"><li>■ Mathieu Jan (INRIA/Région Bretagne)</li><li>■ Sébastien Monnet (ACI MD - GDS)</li><li>■ Loic Cudennec (Sun Microsystems/ Région Bretagne)</li></ul>	<ul style="list-style-type: none"><li>■ Jean-François Deverge (Master 2, IFSIC)</li><li>■ David Noblet (UNH, USA)</li><li>■ Chester Tse (MIT, USA)</li><li>■ Arvind Saraf (MIT, USA)</li></ul>



# Projet REGAL

INRIA Rocquencourt - LIP6

- Domaine de recherche :
  - Adaptation des systèmes aux environnements à large échelle
  - Architectures cibles : GRID / P2P
  
- Compétences de REGAL (relatives à GDS) :
  - Algorithmique répartie
  - Tolérance aux fautes (observation et détection)
  - Réplication du code et des données
  
- Personnels impliqués directement dans GDS :

<b>Permanents</b>	<b>Doctorants</b>	<b>Stagiaire</b>
L. Arantes (MDC Paris 6) P. Sens (Prof. Paris 6)	Marin Bertier (MESR)* Jean-Michel Busca (INRIA) Fabio Picconi (ACI MD - GDS) Julien Sopena (MESR)	Daniel Myers (Fullbright, USA)

\* Actuellement MdC INSA de Rennes, projet PARIS



# Etat du projet à T0+24

---

- 10 réunions GDS entre septembre 2003 et novembre 2005
- Visites des doctorants (sur 2-3 jours)
  - 3 visites de Mathieu Jan (PARIS) chez GRAAL
  - 2 visite de Fabio Picconi (REGAL) chez PARIS
  - 2 visites de Sébastien Monnet (PARIS) chez REGAL
- Logiciels en cours de développement :
  - PARIS : plate-forme JuxMem, outil de déploiement pair-à-pair JDF
  - REGAL : simulateur de systèmes à large échelle LS3, système de fichiers P2P Pastis
  - GRAAL : plate-forme DIET, outil de déploiement GoDIET et de visualisation VizDIET
- Résultats obtenus
  - Détecteur hiérarchique de fautes (REGAL) intégré dans JuxMem (PARIS)
  - Expérimentations préliminaires DIET (GRAAL)/JuxMem (PARIS)
- Travaux en cours
  - Premier prototype DIET/JuxMem, validation à l'aide l'application TLSE
  - Tests d'extensibilité : protocoles de cohérence tolérants aux fautes, algorithmes d'allocation d'espace de stockage



# Etat du projet à T0+24 (suite)

---

- Plate-forme matérielle

- GRAAL : 13 nœuds + plate-forme Grid'5000 Lyon (112 nœuds)
- PARIS : nœuds GDS intégrés dans la plate-forme Grid'5000 (264 nœuds bi-processeurs)
- REGAL : Mise en place d'une plate-forme d'émulation de système à large échelle : 20 machines (22 nœuds bi-processeurs de calcul, 18 mono-processeurs avec 4 ralentisseurs de réseau dummynet)

- Visites extérieures

- Visite de G. Antoniu chez *Sun Microsystems* à Santa Clara (équipe JXTA) – novembre 2003
- Séjours de Mathieu Jan (3 mois) et de Loïc Cudennec (2 mois) chez *Sun Microsystems* - 2005
- Visite de F. Picconi à *l'Université de Rice* (équipe de P. Drushel) – juillet 2004
- Séjour de F. Picconi à *l'Université de Rutgers* (équipe de L. Iftode) – Janvier-Février 2005

- Invités extérieurs

- Stage de 9 semaines de David Noblet (*Université du New Hampshire*) chez PARIS – été 2004
- Visite de P. Hatcher (*Université du New Hampshire*) chez PARIS – mai 2004
- Visite d'André Schiper (*EPFL*) chez REGAL – juin 2004



# Publications GDS à T0+24

---

- Un chapitre de livre à paraître (*Future Generation Grids*, Springer)
- Journaux internationaux (2)
  - 2005 : Kluwer Journal of Supercomputing, Scalable Computing: Practice and Experience, Concurrency and Computation : Practice and Experience
  - 2005 : Journal of Parallel and Distributed Computing
  - 2005 : International Journal of High Performance Computing Applications
  - Soumission : Journal of Grid Computing
- Conférences internationales (8)
  - 2003 : PACT (WIP), IPDPS
  - 2004 : ICPADS, CCGRID, SBAC, Euro-Par
  - 2005 : HPCC, Europar (2), e-Science Grid
- Colloques internationaux (6)
  - 2003 : AGRIDM (PACT)
  - 2004 : AGRIDM (PACT), HCW, GP2PC (CC-GRID), WGAPI (GGF-12), WGC
- Communications nationales (5)
  - 2003 : RenPar
  - 2004 : Ecole DRUIDE (3), GridUse
  - 2005 : CFSE, Journées CDUR





# Collaborations nationales

---

- ACI GRID DataGraal (animation)
  - Ecolé thématique DRUIDE 2004 : Distribution de données à grande échelle (CNRS, INRIA, GDR ARP, Univ. Rennes 1), Le Croisic, mai 2004
- ACI MD
  - GdX – expérimentations GDS/GdX prévues pour 2005/2006 (émulation)
  - MDP2P - workshop sur la gestion de données P2P (mars 2005)
  - Gédéon (protocole de réplication)
- AS : Algorithmique Distribuée et Applications
  - Journées thématiques « algorithmique distribuée et applications », Porquerolles, septembre 2004
- Grid'5000
  - Déploiement et tests à grande échelle multi-sites



# Collaborations internationales

---

- Partenaires académiques
  - AIST, Tokyo (partage transparent P2P à grande échelle)
  - Réseau Européen d'Excellence CoreGRID / Université de Pise (stockage de données sur grille)
  - Univ. Rice (système P2P)
  - Univ. Libre d'Amsterdam (réplication adaptable)
  - EPFL (détection de fautes)
- Partenaires industrielles
  - Sun Microsystems : JuxMem/JXTA/ Grid'5000
    - Financement d'une thèse à l'IRISA démarrée en 2005 (Loïc Cudennec)



# Pour plus de détails...

---

- Site web : <http://www.irisa.fr/GDS/>
  - Programme des réunions
  - Présentations en ligne
  - Visites
  - Publications
  - Archive liste de diffusion : [acimd-gds@irisa.fr](mailto:acimd-gds@irisa.fr)